Docket No.: POU920030100US1

Inventor: Joseph F. Skovira

Title: BALANCING WORKLOAD
OF A GRID COMPUTING
ENVIRONMENT

APPLICATION FOR UNITED STATES

LETTERS PATENT

"Express Mail" Mailing Label No.: ER 046554978 US

Date of Deposit:

August 5, 2003

I hereby certify that this paper is being deposited with the United States Postal Service as "Express Mail Post Office to Addressee" service under 37 CFR 1.10 on the date indicated above and is addressed to: Mail Stop: Patent Application, Commissioner for Patents, Alexandria, VA 22313-1450

Name: SUSAN L. PHELPS

Signature

INTERNATIONAL BUSINESS MACHINES CORPORATION

BALANCING WORKLOAD OF A GRID COMPUTING ENVIRONMENT

Technical Field

[0001] This invention relates, in general, to grid computing, and in particular, to managing the workload of a grid computing environment.

Background of the Invention

[0002] A grid computing environment allows the interconnection of a plurality of heterogeneous and/or geographically distant systems. To facilitate interconnecting the systems, in one example, a Globus toolkit offered by International Business Machines Corporation, Armonk, New York, is employed. Globus allows a user to specify which system of the plurality of systems is to run a job. The user submits a job to the selected system using a Resource Specification Language (RSL). In response to Globus receiving the RSL, Globus converts the RSL into a correct format for the scheduler on the target system. For example, if the scheduler is LoadLeveler offered by International Business Machines Corporation, then the RSL is converted into a command file.

[0003] Since users can select the one or more systems to run their jobs or despite the selections, the systems of a grid computing environment can become unbalanced. For instance, one system may have too much work, while another system may have too little. Thus, a need exists for a capability to balance the workload of a grid computing environment. A further need exists for a capability for determining the best fit for particular work.

Summary of the Invention

[0004] The shortcomings of the prior art are overcome and additional advantages are provided through the provision of a method of balancing workload of a computing environment. The method includes, for instance, obtaining information regarding one or more systems of a plurality of systems of a grid computing environment; and balancing

workload of at least two systems of the plurality of systems using at least a portion of the obtained information.

[0005] System and computer program products corresponding to the abovesummarized method are also described and claimed herein.

[0006] Additional features and advantages are realized through the techniques of the present invention. Other embodiments and aspects of the invention are described in detail herein and are considered a part of the claimed invention.

Brief Description of the Drawings

[0007] The subject matter which is regarded as the invention is particularly pointed out and distinctly claimed in the claims at the conclusion of the specification. The foregoing and other objects, features, and advantages of the invention are apparent from the following detailed description taken in conjunction with the accompanying drawings in which:

- [0008] FIG. 1 depicts one embodiment of a computing environment incorporating and using one or more aspects of the present invention;
- [0009] FIG. 2 depicts one embodiment of the logic associated with balancing workload of the computing environment of FIG. 1, in accordance with an aspect of the present invention;
- [0010] FIG. 3 depicts further details regarding one embodiment of the logic associated with workload balancing, in accordance with an aspect of the present invention; and
- [0011] FIG. 4 depicts one embodiment of the logic used to determine which system of the environment is to run a given job, in accordance with an aspect of the present invention.

Best Mode for Carrying Out the Invention

[0012] In accordance with an aspect of the present invention, workload balancing is performed in a grid computing environment. In one example, a manager daemon of the grid computing environment obtains information regarding one or more systems of the environment, and determines based on the obtained information placement of workload on those systems. The placement of workload can include, for instance, moving a job from one system to another or initially placing a job on a particular system, etc. As one example, the information is obtained from schedulers of the systems.

[0013] Grid computing enables the virtualization of distributed computing and data resources such as processing, network bandwidth and storage capacity to create a single system image, granting users and applications seamless access to vast information technology (IT) capabilities. Often, the systems of a grid computing environment are heterogeneous systems. That is, at least one system of the plurality of systems of the environment includes different hardware and/or software from at least one other system of the environment. Additionally or alternatively, the systems may be geographically distant from one another. Further details regarding grid computing may be found, for instance, at www-1.ibm.com/grid/about_grid/what_is.shtml.

[0014] One embodiment of a grid computing environment incorporating and using one or more aspects of the present invention is depicted in FIG. 1. A grid computing environment 100 includes, for instance, a plurality of systems 102. In this particular example, two systems, System A and System B, are depicted. However, in other examples, more than two systems are included in the computing environment. In one example, System A includes a Scalable Parallel (SP) machine having a plurality of RS/6000 nodes, which is offered by International Business Machines Corporation, Armonk, New York, and System B includes a LINUX cluster also offered by International Business Machines Corporation. Systems 102 are coupled to one another via a connection 104, such as, for instance, an Ethernet connection or other type of connection.

[0015] System 102 includes, for instance, a scheduler 106 used in scheduling jobs on the system. A scheduler can be one of many types of schedulers and each system may have the same type of scheduler or different type of scheduler. As one example, scheduler 106 of System A includes LoadLeveler offered by International Business Machines Corporation, and scheduler 106 of System B includes the Portable Batch System (PBS) offered by Altair Grid Technologies, LLC. One example of LoadLeveler is described in an IBM publication entitled, "IBM LoadLeveler: Using and Administering," V3R1, IBM Pub. No. SA22-7881-00, December 2001, which is hereby incorporated herein by reference in its entirety.

[0016] In one example, at least one scheduler performs backfill scheduling. Backfill scheduling allows an application to run out of order as long as it does not affect the start time of an application already scheduled to execute. One example of backfill scheduling is described in U.S. Patent Application No. 10/406,985 entitled "Backfill Scheduling Of Applications Based On Data Of The Applications," filed April 4, 2003, which is hereby incorporated herein by reference in its entirety.

[0017] Since, in one example, the systems of the grid computing environment are heterogeneous, a toolkit, offered by International Business Machines Corporation and referred to as Globus, is used to facilitate communication between the systems. The toolkit creates a common layer between the systems. For instance, for a Globus enabled system, information for a job passes through Globus, Globus converts it to a Globus format, passes the information to another Globus system, which then converts the information into a form known by the receiving system. This allows systems having one or more of different operating systems, different middleware and/or different schedulers to effectively communicate. Further details regarding Globus can be found, for instance, in "Enabling Applications for Grid Computing with Globus," IBM publication no. SG24-6936-00, June 18, 2003, which is hereby incorporated herein by reference in its entirety.

[0018] In accordance with an aspect of the present invention, one of the systems in the grid computing environment also includes a manager daemon 108. The manager

daemon runs in the background and is responsible for balancing the workload among at least a portion of the systems of the environment. The manager daemon obtains (e.g., is provided, determines, etc.) information regarding a plurality of systems to be managed. This information includes, for instance, identification of the systems, manner in which to contact the systems, etc.

[0019] The manager daemon periodically executes logic to balance the workload of a grid computing environment. In one example, this logic is executed at configurable time intervals (e.g., every 5 minutes). As another example, execution of the logic is event based (e.g., upon startup and/or completion of a job, change in available system resources, etc.). One embodiment of the logic associated with balancing the workload of a grid computing environment is described with reference to FIGs. 2-4.

[0020] Referring initially to FIG. 2, the manager daemon obtains scheduler information for one or more systems, STEP 200. For instance, the manager daemon contacts the schedulers of those systems to obtain desired information. This information includes, for instance, the current free nodes of the system, the job queue of waiting jobs for that system, and scheduler specific variable settings for the current state of the system job mix, such as a shadow time for the next waiting job (i.e., how long does the job need to wait for resources) and one or more resources protected by the shadow time.

[0021] Based on the information obtained, the manager daemon performs workload balancing, STEP 202. Further details regarding one example of workload balancing is described with reference to FIG. 3. Initially, the scheduling information is used to determine which system is to run a given job, STEP 300. In one example, this includes determining which idle jobs on a particular system might run on another system. One example of the logic employed to make this determination is described with reference to FIG. 4. In the example described herein, a determination is made as to whether one or more jobs on System A can be moved to System B. However, it will be apparent to those skilled in the art that similar logic is used to move jobs to System A or to other systems being managed.

[0022] Referring to FIG. 4, a determination is made as to whether there are any free nodes on System B, INQUIRY 400. If there are no free nodes, then processing is complete, STEP 402. However, if there are one or more free nodes, then a further determination is made as to whether there are one or more idle jobs on System A, INQUIRY 404. Should there be an idle job on System A, then a further determination is made as to whether the idle job fits on System B, INQUIRY 406. If the idle job does fit on System B, then, in one example, a further determination is made as to whether the job can backfill, INQUIRY 408. If the job does fit on the new system and can backfill, then the job is placed on a transfer list, STEP 410. Otherwise, a determination is made as whether there are further idle jobs on System A, INQUIRY 404. If not, then processing is complete, STEP 402.

[0023] Returning to FIG. 3, in addition to determining which system is to run a given job, workload balancing further includes placing the job on that system, STEP 302. In one example, this includes moving each job (or a portion of the jobs) from the transfer list to the indicated system(s). This includes, for instance, placing the job on hold in the original system (e.g., System A) to prevent the job selected for transfer from starting. The job is then submitted to the new system (e.g., System B). If the move is successful, then the job is cancelled from the first system. By using a technique of hold and then move, further error checking may be provided at the discretion of the designer. In one example, commands provided by Globus are used in the moving.

[0024] Described in detail above is one embodiment of the logic associated with using a daemon to perform workload balancing in a grid computing environment. One embodiment of pseudo-code used to perform this workload balancing is presented below:

Do forever {

Get a current snapshot of the 2 batch systems
Access LoadLeveler on system A for FreeNodesA, ShadowTimeA, IdleJobsA
Access LoadLeveler on system B for FreeNodesB, ShadowTimeB, IdleJobsB
Clear the Transfer Lists A2B and B2A

Find out which Idle jobs on system A might run on system B

```
if (FreeNodesB) { # if there are any free nodes in system B...
   Foreach (IdleJobsA) { # Then for all the idle jobs on system A...
     If (JobA node requirement <= FreeNodesB) { # if the job fits in system B...
       If (JobA Wallclock time <= ShadowTimeB) { # if the job can backfill...
         Place JobA on the Transfer List A2B
     }
   }
 # Find out which Idle jobs on system B might run on system A
 if (FreeNodesA) { # if there are any free nodes in system A...
   Foreach (IdleJobsB) { # Then for all the idle jobs on system B...
     If (JobB node requirement <= FreeNodesA) { # if the job fits in system A...
       If (JobB Wallclock time <= ShadowTimeA) { # if the job can backfill...
         Place JobB on the Transfer List B2A
   }
 # Move potential jobs from A to B
 foreach (job in the A2B array) {
   Move JobA to SystemB
 # Move potential jobs from B to A
 foreach (job in the B2A array) {
   Move JobB to SystemA
 Sleep for a short time # User configurable, about 30 seconds
} # end of Do forever
# Move Job subroutine which moves a job from one system to another
sub Move JobX to SystemY {
 Place JobX on System Hold
 Submit JobX to SystemY
 Once JobX appears on SystemY {
   Remove JobX from SystemX
} # end of subroutine
```

[0025] Described herein is a capability for balancing the workload of a grid computing environment. To balance the workload, in one example, work is moved from one system that is more heavily loaded to another system that is more lightly loaded. In other examples, workloads are balanced in other ways. For example, workload balancing

may include initially determining which system is to run a particular job and submitting the job on that system. In that case, users submit jobs into a holding pen, which is visible to the daemon. Although the jobs in the holding pen are visible to the daemon, in this example, they are not visible to the schedulers on the individual systems. The daemon requests information from the schedulers and based on this information determines the best fit for a particular job. The daemon then submits the job to the selected system.

[0026] Although the initial submission of jobs is controlled, the systems can still become unbalanced. This imbalance may occur because of unpredictable events during job execution (e.g., job failure which causes the job to complete earlier than expected), which would disrupt previous queuing decisions, etc. Thus, the daemon also executes the above logic, in one example, to maintain workload balance.

[0027] The information used in balancing the workload can be different, less and/or in addition to that described above. As examples, job class and/or resource matches (such as memory or software licenses), as well as other information could be used to decide the placement of jobs.

[0028] The workload balancing capability of the present invention advantageously enables workload on two or more systems of a grid computing environment to be balanced. Again, although two systems are described herein, more than two systems with independent batch queuing capabilities could be controlled with a single daemon. The logic would be expanded to examine information from the additional systems. Further, although examples of systems are provided above, many other possibilities exist. As one example, the systems are homogeneous, but are geographically distant. Many other variations also exist.

[0029] In one aspect, the daemon can be deactivated. When it is deactivated, the users can still submit jobs on a plurality of systems, but automatic load balancing between two grid connected systems does not occur.

[0030] Further, although the above example uses a backfill scheduling technique, other scheduling techniques including those that do not backfill may be used. If a technique that does not use backfill is employed, then shadow time may not be included in the collected information. For example, in a FIFO scheduling technique, the daemon determines idle nodes, free jobs, and possibly idle job order, but it does not require shadow time. When deciding to move jobs to a system, the free resources are considered and there is no shadow time test. In a like manner, other batch scheduling techniques can be used in managing workloads.

[0031] Further, for those schedulers using backfill, in another embodiment, the list of which resources are protected by shadow time (and which are not) is employed to improve the decision making process. For example, a job can be transferred which has a wallclock estimate greater than the shadow time to nodes which are not protected by shadow time (hence, not limited to backfill timing constraints).

[0032] Moreover, although examples of schedulers are provided above, many other schedulers can be used without departing from the spirit of the present invention.

Examples of other schedulers include, for instance, the Load Sharing Facility (LSF) offered by Platform Computing, and Maui, offered by the Maui Supercomputing Center.

[0033] As a further embodiment, more than one system may include a manager daemon. One may be backup of another and/or multiple daemons may work together to manage the workload of the grid computing environment, etc. Moreover, one or more systems of the environment may not have a scheduler, but instead, are scheduled by other schedulers, etc.

[0034] Advantageously, one or more aspects of the present invention enable a grid computing environment to be workload balanced. This increases efficiency and productivity. By being dynamic and automatic, the balancing is transparent to users. By obtaining the information from schedulers and maintaining the scheduling responsibilities with the schedulers, complexity of the manager daemon is minimized. Since the

information obtained by the daemon comes from complex scheduling software, the amount of information input to the daemon is reduced. Further, the scheduler can send the results of an already executed algorithm to the daemon and the daemon does not have to perform the complex analysis (such as computation of shadow time, etc.).

[0035] Advantageously, one or more aspects of the present invention enable a plurality of parallel machines, in which each machine is independently administered, to combine resources under, for instance, a single Globus implementation.

[0036] The present invention can be included in an article of manufacture (e.g., one or more computer program products) having, for instance, computer usable media. The media has therein, for instance, computer readable program code means or logic (e.g., instructions, code, commands, etc.) to provide and facilitate the capabilities of the present invention. The article of manufacture can be included as a part of a computer system or sold separately.

[0037] Additionally, at least one program storage device readable by a machine embodying at least one program of instructions executable by the machine to perform the capabilities of the present invention can be provided.

[0038] The flow diagrams depicted herein are just examples. There may be many variations to these diagrams or the steps (or operations) described therein without departing from the spirit of the invention. For instance, the steps may be performed in a differing order, or steps may be added, deleted or modified. All of these variations are considered a part of the claimed invention.

[0039] Although preferred embodiments have been depicted and described in detail herein, it will be apparent to those skilled in the relevant art that various modifications, additions, substitutions and the like can be made without departing from the spirit of the invention and these are therefore considered to be within the scope of the invention as defined in the following claims.